

PromptBERT: Improving BERT Sentence Embeddings with Prompts

Anonymous ACL submission

Abstract

The poor performance of the original BERT¹ for sentence semantic similarity has been widely discussed in previous works. We find that unsatisfactory performance is mainly due to the static token embeddings biases and the ineffective BERT layers, rather than the high cosine similarity of the sentence embeddings.

To this end, we propose a prompt based sentence embeddings method which can reduce token embeddings biases and make the original BERT layers more effectively. By reformulating the sentence embeddings task as the fill-in-the-blanks problem, our method significantly improves the performance of original BERT. We discuss two prompt representing methods and three prompt searching methods for prompt based sentence embeddings. Moreover, we propose a novel unsupervised training objective by the technology of template denoising, which substantially shortens the performance gap between the supervised and unsupervised setting. For experiments, we evaluate our method on both non fine-tuned and fine-tuned settings. Even a non fine-tuned method can outperform the fine-tuned methods like unsupervised ConSERT on STS tasks. Our fine-tuned method outperforms the state-of-the-art method SimCSE in both unsupervised and supervised settings. Compared to SimCSE, we achieve 2.29 and 2.58 points improvements on BERT and RoBERTa respectively under the unsupervised setting.

1 Introduction

In recent years, we have witnessed the success of pre-trained language models like BERT (Devlin et al., 2018) and RoBERTa (Liu et al., 2019) in sentence embeddings (Gao et al., 2021; Yan et al., 2021). However, the original BERT still shows poor performance in sentence embeddings

¹In this paper, we use “original BERT” to refer to the BERT like models, which are not fine-tuned on downstream tasks.

(Reimers and Gurevych, 2019; Li et al., 2020). The most commonly used example is that it underperforms the traditional word embedding methods like GloVe (Pennington et al., 2014).

Previous research has linked anisotropy to explain the poor performance of the original BERT (Li et al., 2020; Yan et al., 2021; Gao et al., 2021). Anisotropy makes the token embeddings occupy a narrow cone, resulting in a high similarity between any sentence pair (Li et al., 2020). Li et al. (2020) proposed a normalizing flows method to transform the sentence embeddings distribution to a smooth and isotropic Gaussian distribution and Yan et al. (2021) presented a contrastive framework to transfer sentence representation. The goal of these methods is to eliminate anisotropy in sentence embeddings. However, we find that anisotropy is not the primary cause of poor semantic similarity. For example, averaging the last layer of the original BERT is even worse than averaging its static token embeddings in semantic textual similarity task, but the sentence embeddings from last layer are less anisotropic than static token embeddings.

Following this result, we find the original BERT layers² actually damage the quality of sentence embeddings. However, if we treat static token embeddings³ as word embedding, it still yields unsatisfactory results compared to GloVe. Inspired by (Li et al., 2020), who found token frequency biases its distribution, we find the distribution is not only biased by frequency, but also case sensitive and subword in WordPiece (Wu et al., 2016). We design a simple experiment to test our conjecture by simply removing these biased tokens (e.g., high frequency subwords and punctuation) and using the average of the remaining token embeddings as sentence representation. It can outperform the

²We denote the transformer blocks in BERT as BERT layers.

³We denote the BERT token embeddings as static token embeddings.

077 Glove and even achieve results comparable to post-
078 processing methods BERT-flow (Li et al., 2020)
079 and BERT-whitening (Su et al., 2021).

080 Motivated by these findings, avoiding embedding
081 bias can improve the performance of sentence
082 representations. However, it is labor-intensive to
083 manually remove embedding biases and it may re-
084 sult in the omission of some meaningful words
085 if the sentence is too short. **Inspired by** (Brown
086 et al., 2020), which has reformulated the different
087 NLP tasks as fill-in-the-blanks problems by differ-
088 ent prompt, we propose a prompt based method by
089 using the template to obtain the sentence repres-
090 entations in BERT. Prompt based method can avoid
091 embedding bias and utilize the original BERT lay-
092 ers. We find original BERT can achieve reasonable
093 performance with the help of the template in sen-
094 tence embeddings, and it even outperforms some
095 BERT based methods, which fine-tune BERT in
096 down-stream tasks.

097 Our approach is equally applicable to fine-tuned
098 setting. Current methods utilize the contrastive
099 learning to help the BERT learn better sentence em-
100 beddings (Gao et al., 2021; Yan et al., 2021). **How-**
101 **ever, the unsupervised methods still suffer from**
102 **leaking proper positive pairs.** Yan et al. (2021)
103 discuss four data augmentation methods, but the
104 performance seems worse than directly using the
105 dropout in BERT as noise (Gao et al., 2021). We
106 find the prompts can provide a better way to gener-
107 ate positive pairs by different viewpoints from dif-
108 ferent templates. To this end, we propose a prompt
109 based contrastive learning method with template
110 denoising to leverage the power of BERT in an
111 unsupervised setting, which significantly shortens
112 the gap between the supervised and unsupervised
113 performance. Our method achieves the state-of-
114 the-art results in both unsupervised and supervised
115 settings.

116 2 Related Work

117 Learning sentence embeddings as a fundamental
118 NLP problem has been largely studied. Currently,
119 how to leverage the power of BERT in sentence em-
120 beddings has become a new trend. Many works (Li
121 et al., 2020; Gao et al., 2021) achieved strong per-
122 formance with BERT in both supervised and unsup-
123ervised settings. Among these works, contrastive
124 learning based methods achieve the state-of-the-art
125 results. These works (Gao et al., 2021; Yan et al.,
126 2021) pay attention to constructing positive sen-

tence pairs. Gao et al. (2021) proposed a novel
127 contrastive training objective to directly use inner
128 dropout as noise to construct positive pairs. Yan
129 et al. (2021) discuss four methods to construct pos-
130 itive pairs.

131 Although BERT achieved great success in sen-
132 tence embeddings, the original BERT shows unsat-
133 isfactory performance. Contextual token embed-
134 dings from original BERT even underperform the
135 word embeddings like GloVe. One explanation is
136 the anisotropy in the original BERT, which causes
137 sentence pairs to have high similarity. Following
138 this explanation, BERT-flow (Li et al., 2020) and
139 BERT-whitening (Su et al., 2021) have been pro-
140 posed to reduce the anisotropy by post-processing
141 the sentence embeddings from original BERT.

143 3 Rethinking the Sentence Embeddings of 144 Original BERT

145 Previous works (Yan et al., 2021; Gao et al., 2021)
146 explained the poor performance of original BERT
147 is limited by the learned anisotropic token embed-
148 dings space, where the token embeddings occupy
149 a narrow cone. However, we find that anisotropy
150 is not a key factor to inducing poor semantic sim-
151 ilarity by examining the relationship between the
152 anisotropy and performance. We think the main
153 reasons are the ineffective BERT layers and static
154 token embedding biases.

155 **Observation 1: Original BERT layers fail to**
156 **improve the performance.** In this section, we an-
157alyze the influence of BERT layers by comparing
158 the two sentence embeddings methods: averaging
159 static token embeddings (input of the BERT lay-
160 ers) and averaging last layer (output of the BERT
161 layers). We report the sentence embeddings perfor-
162 mance and its sentence level anisotropy.

163 To measure the anisotropy, we follow the work
164 of (Ethayarajh, 2019) to measure the sentence level
165 anisotropy in sentence embeddings. Let s_i be a
166 sentence that appears in corpus $\{s_1, \dots, s_n\}$. The
167 anisotropy can be measured as follows:

$$\frac{1}{n^2 - n} \left| \sum_i \sum_{j \neq i} \cos(M(s_i), M(s_j)) \right| \quad (1)$$

168 where M denotes the sentence embeddings method,
169 which maps the raw sentence to its embedding
170 and \cos is the cosine similarity. In other words,
171 the anisotropy of M is measured by the average
172 cosine similarity of a set of sentences. If sen-

tence embeddings was isotropic (i.e., directionally uniform), then the average cosine similarity between uniformly randomly sampled sentences would be 0 (Arora et al., 2016). The closer it is to 1, the more anisotropic the embedding of sentences. We randomly sample 100,000 sentences from the Wikipedia corpus to compute the anisotropy.

We compare different pre-trained models (*bert-base-uncased*, *bert-base-cased* and *roberta-base*) and different sentence embeddings methods (last layer average, averaging of last hidden layer tokens as sentence embeddings and static token embeddings, directly averaging of static token embeddings). We have shown the spearman correlation, sentence level anisotropy of these methods in Table 1.

Pre-trained models	Correlation	Sentence anisotropy
<i>Static token embeddings avg.</i>		
<i>bert-base-uncased</i>	56.02	0.8250
<i>bert-base-cased</i>	56.65	0.5755
<i>roberta-base</i>	55.88	0.5693
<i>Last layer avg.</i>		
<i>bert-base-uncased</i>	52.57	0.4874
<i>bert-base-cased</i>	56.93	0.7514
<i>roberta-base</i>	53.49	0.9554

Table 1: The spearman correlation, sentence anisotropy of Last layer average. and Static token embeddings average. The spearman correlation is the average of correlation on STS12-16, STS-B and SICK.

As Table 1 shows, we find the BERT layers in *bert-base-uncased* and *roberta-base* significantly harm the sentence embeddings performance. Even in *bert-base-cased*, the gain of BERT layers is trivial with only 0.28 improvement. We also show the sentence level anisotropy of each method. The performance degradation of the BERT layers seems not to be related to the sentence level anisotropy. For example, the last layer average is more isotropic than the static token embeddings average in *bert-base-uncased*. However, the static token embeddings average achieves better sentence embeddings performance.

Observation 2: Embedding biases harms the sentence embeddings performance. Li et al. (2020) found that token embeddings can be biased to token frequency. Similar problems have been studied in (Yan et al., 2021). The anisotropy in BERT static token embeddings is sensitive to token frequency. Therefore, we investigate whether

embedding bias yields unsatisfactory performance of sentence embeddings. We observe that the token embeddings is not only biased by token frequency, but also subwords in WordPiece (Wu et al., 2016) and case sensitive.

As shown in Figure 1, we visualize these biases in the token embeddings of *bert-base-uncased*, *bert-base-cased* and *roberta-base*. The token embeddings of three pre-trained models are highly biased by the token frequency, subword and case. The token embeddings can roughly divided into three regions according to the subword and case biases : 1) the lowercase begin-of-word tokens, 2) the uppercase begin-of-word tokens and 3) the subword tokens. For uncased pre-trained model *bert-base-uncased*, the token embeddings also can roughly divided into two regions: 1) the begin-of-word tokens, 2) the subword tokens.

For frequency bias, we can observe that high frequency tokens are clustered, while low frequency tokens are dispersed sparsely in all models (Yan et al., 2021). The begin-of-word tokens are more vulnerable to frequency than subword tokens in BERT. However, the subword tokens are more vulnerable in RoBERTa.

Previous works (Yan et al., 2021; Li et al., 2020) often connect the concept of "token embeddings bias" with the token embeddings anisotropy as the reason for bias. However, we think the anisotropy is unrelated to the bias. The bias means the distribution of embedding is disturbed by some irrelevant information like token frequency, which can be directly visualized according to the PCA. For the anisotropy, it means the whole embedding occupies a narrow cone in the high dimensional vector space, which cannot be directly visualized.

<i>M</i>	average cosine similarity
<i>bert-base-uncased</i>	0.4445
<i>bert-base-cased</i>	0.1465
<i>roberta-base</i>	0.0235

Table 2: The average cosine similarity in static token embeddings

Table 2 shows the static token embeddings anisotropy of three pre-trained models in Figure 1 according to the average the cosine similarity between any two token embeddings. Contrary to the previous conclusion (Yan et al., 2021; Li et al., 2020), we find only *bert-base-uncased*'s static token embeddings is highly anisotropic. The static

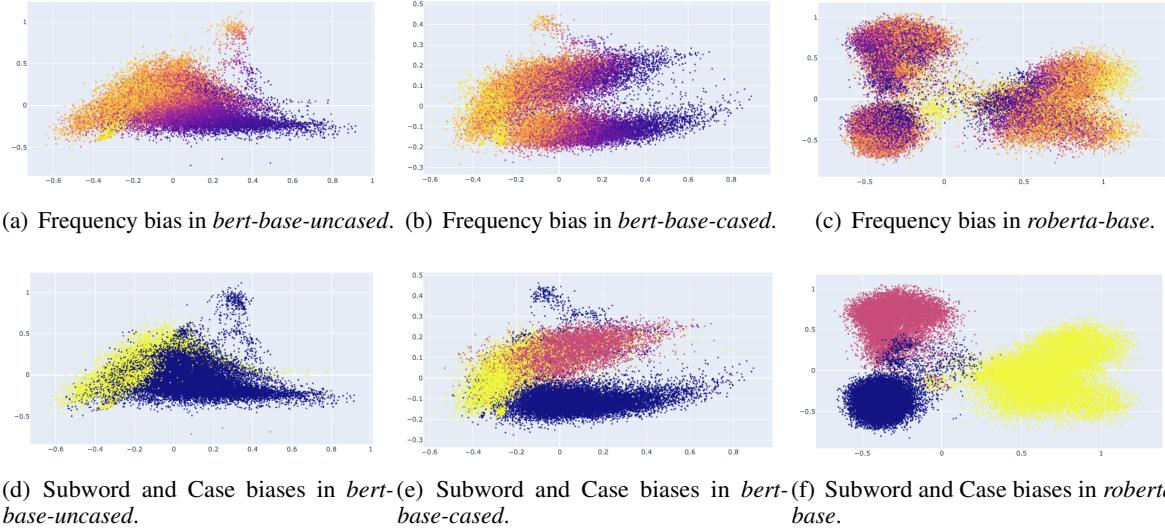


Figure 1: 2D visualization of token embeddings with different biases. For frequency bias, the darker the color, the higher the token frequency. For subword and case bias, yellow represents subword and red represents the token contains capital letters.

253 token embeddings like *roberta-base* are isotropic
254 with 0.0235 average cosine similarity. For bi-
255 ases, these models are suffered from the biases
256 in static token embeddings, which is irrelevant to
257 the anisotropy.

270 Manually removing embedding biases is a sim-
271 ple method to improve the performance of sentence
272 embeddings. However, if the sentence is too short,
273 this is not an adequate solution, which may result
274 in the omission of some meaningful words.

	cased	uncased	roberta
Static Token Embeddings	56.93	56.02	55.88
– Freq.	60.27	59.65	65.41
– Freq. & Sub.	64.83	62.20	64.89
– Freq. & Sub. & Case	65.07	-	65.06
– Freq. & Sub. & Case & Pun.	66.05	63.10	67.64

Table 3: The influence of static embedding biases in spearman correlation. The spearman correlation is the average of STS12-16, STS-B and SICK. Cased, uncased and roberta represent *bert-base-cased*, *bert-base-uncased* and *roberta-base*. For Freq., Sub., Case. and Pun., we remove the top frequency tokens, subword tokens, uppercase tokens and punctuation respectively. More details can be found in Appendix A.

258 To prove the negative impact of biases, we show
259 the influence of biases to the sentence embeddings
260 with averaging static token embeddings as sentence
261 embeddings (without BERT layers). The results of
262 eliminating embedding biases are quite impressive
263 on three pre-trained models in Table 3. Simply
264 removing a set of tokens, the result can be im-
265 proved by 9.22, 7.08 and 11.76 respectively. The
266 final result of *roberta-base* can outperform post-
267 processing methods such as BERT-flow (Li et al.,
268 2020) and BERT-whitening (Su et al., 2021) with
269 only using static token embeddings.

4 Prompt Based Sentence Embeddings

Inspired by (Brown et al., 2020), we propose a prompt based sentence method to obtain sentence embeddings. By reformulating the sentence embedding task as the mask language task, we can effectively use the original BERT layers by leveraging the large-scale knowledge. We also avoid the embedding biases by representing sentences from [MASK] tokens.

However, unlike the text classification or question-answering tasks, the output in sentence embeddings is not the label tokens predicted by MLM classification head, but the vector to represent the sentence. We follow these two problems to discuss the implementation of prompt based sentence embeddings: 1) how to represent sentences with the prompt, and 2) how to find a proper prompt for sentence embeddings. Based on these, we propose a prompt based contrastive learning method to fine-tuning BERT on sentence embeddings.

4.1 Represent Sentence with the Prompt

In this section, we discuss two methods to represent one sentence with a prompt. For example, we have a template “[X] means [MASK]”, where [X] is a

placeholder to put sentences and [MASK] represent the [MASK] token. Given a sentence x_{in} , we map x_{in} to x_{prompt} with the template. Then we feed x_{prompt} to a pre-trained model to generate sentence representation \mathbf{h} .

One method is to use the hidden vector of [MASK] token as sentence representation:

$$\mathbf{h} = \mathbf{h}_{[MASK]} \quad (2)$$

For the second method like other prompt based tasks, we get the top- k tokens according to $\mathbf{h}_{[MASK]}$ and MLM classification head, then find the weighted average of these tokens according to probability distribution. The \mathbf{h} can be formulated as:

$$\mathbf{h} = \frac{\sum_{v \in \mathcal{V}_{top-k}} \mathbf{W}_v P([MASK] = v | \mathbf{h}_{[MASK]})}{\sum_{v \in \mathcal{V}_{top-k}} P([MASK] = v | \mathbf{h}_{[MASK]})} \quad (3)$$

where v is the BERT token in the top- k tokens set \mathcal{V}_{top-k} , \mathbf{W}_v is the static token embeddings of v and $P([MASK] = v | \mathbf{h}_{[MASK]})$ denotes the probability of token v be predicted by MLM head with $\mathbf{h}_{[MASK]}$.

The second method, which maps the sentence to the tokens, is more conventional than the first. But its disadvantages are obvious: 1) as previously noted, due to the sentence embeddings from averaging of static token embeddings, it still suffers from biases. 2) weight averaging makes the BERT hard to fine-tune in down-stream tasks. For these reasons, we represent the sentence with the prompt by the first method.

4.2 Prompt Search

For prompt based tasks, one key challenge is to find templates. We discuss three methods to search template in this section: manual search, template generation based T5 (Gao et al., 2020) and OptiPrompt (Zhong et al., 2021). We use the spearman correlation in the STS-B development set as the main metric to evaluate different templates.

For manual search, we need to hand-craft templates, which encourage the whole sentence to be represented in $\mathbf{h}_{[MASK]}$. To search templates, we divide the template into two parts: relationship tokens, which denotes the relationship between [X] and [MASK], and prefix tokens, which wraps [X]. Then we greedily search for templates following the relationship tokens and prefix tokens.

Some results of greedy searching are shown in Table 4. When it comes to sentence embeddings,

Template	STS-B dev.
<i>Searching for relationship tokens</i>	
[X] [MASK] .	39.34
[X] is [MASK] .	47.26
[X] mean [MASK] .	53.94
[X] means [MASK] .	63.56
<i>Searching for prefix tokens</i>	
This [X] means [MASK] .	64.19
This sentence of [X] means [MASK] .	68.97
This sentence of “[X]” means [MASK] .	70.19
This sentence : “[X]” means [MASK] .	73.44

Table 4: Greedy searching templates on *bert-base-uncased*.

different templates produce extremely varied results. Compared to simply concatenating the [X] and [MASK], complex templates like *This sentence : “[X]” means [MASK]*, can improve the spearman correlation by 34.10.

For template generation based on T5, Gao et al. (2020) proposed a novel method to automatically generate templates by using T5 to generate templates according to the sentences and corresponding labels. The generated templates can outperform the manual searched templates in the GLUE benchmark (Wang et al., 2018).

However, the main issue to implement it is the lack of label tokens. Tsukagoshi et al. (2021) successfully transformed the sentence embeddings task to the text classification task by classifying the definition sentence to its word according to the dictionary. Inspired by this, we use words and corresponding definitions to generate 500 templates (e.g., orange: a large round juicy citrus fruit with a tough bright reddish-yellow rind). Then we evaluate these templates in the STS-B development set, the best spearman correlation is 64.75 with the template “Also called [MASK]. [X]”. Perhaps it is the gap between sentence embeddings and word definition. This method cannot generate better templates compared to manual searching.

OptiPrompt (Zhong et al., 2021) replaced discrete template with the continuous template. To optimize the continuous template, we use the unsupervised contrastive learning as training objective following the settings in (Gao et al., 2021) with freezing the whole BERT parameters, and the continuous template is initialized by manual template’s static token embeddings. Compared to the input manual template, the continuous template can increase the spearman correlation from 73.44 to 80.90 on STS-B development set.

384
385

4.3 Prompt Based Contrastive Learning with Template Denoising

386 Recently, contrastive learning successfully lever-
387 ages the power of BERT in sentence embeddings.
388 A challenge in sentence embeddings contrastive
389 learning is how to construct proper positive in-
390 stances. Gao et al. (2021) directly used the dropout
391 in the BERT as positive instances. Yan et al. (2021)
392 discussed the four data augmentation strategies
393 such as adversarial attack, token shuffling, cutoff
394 and dropout in the input token embeddings to con-
395 struct positive instances. Motivated by the prompt
396 based sentence embeddings, we propose a novel
397 method to reasonably generate positive instances
398 based on prompt.

399 The idea is using the different templates to repre-
400 sent the same sentence as different points of view,
401 which helps model to produce more reasonable pos-
402 itive pairs. In order to reduce the influence of the
403 template itself on the sentence representation, we
404 propose a novel way to denoise the template infor-
405 mation. Given the sentence x_i , we first calculate
406 the corresponding sentence embeddings \mathbf{h}_i with a
407 template. Then we calculate the template bias $\hat{\mathbf{h}}_i$
408 by directly feeding BERT with the template and
409 the same template position ids. For example, if the
410 x_i has 5 tokens, then the position ids of template
411 tokens after the [X] will be added by 5 to make
412 sure the position ids of template is same. Finally,
413 we can directly use the $\mathbf{h}_i - \hat{\mathbf{h}}_i$ as the denoised sen-
414 tence representation. For the template denoising,
415 more details can be found in Discussion.

416 Formally, let \mathbf{h}'_i and \mathbf{h}_i denote the sentence em-
417 beddings of x_i with different templates, $\hat{\mathbf{h}}'_i$ and $\hat{\mathbf{h}}_i$
418 denotes the two template biases of the x_i respec-
419 tively, the final training objective is as follows:

$$\ell_i = -\log \frac{e^{\cos(\mathbf{h}_i - \hat{\mathbf{h}}_i, \mathbf{h}'_i - \hat{\mathbf{h}}'_i)/\tau}}{\sum_{j=1}^N e^{\cos(\mathbf{h}_i - \hat{\mathbf{h}}_i, \mathbf{h}'_j - \hat{\mathbf{h}}'_j)/\tau}} \quad (4)$$

421 where τ is a temperature hyperparameter in
422 contrastive learning and N is the size of mini-batch.

423

5 Experiments

424 We conduct experiments on STS tasks with non
425 fine-tuned and fine-tuned BERT settings. For non
426 fine-tuned BERT settings, we exploit the per-
427 formance of original BERT in sentence embeddings,
428 which corresponds to the previous findings of the
429 poor performance of original BERT. For fine-tuned

430 BERT settings, we report the unsupervised and
431 supervised results by fine-tuning BERT with down-
432 stream tasks. The results of transfer tasks are in
433 Appendix B.

434

5.1 Dataset

435 Following the past works (Yan et al., 2021; Gao
436 et al., 2021; Reimers and Gurevych, 2019), we con-
437 duct our experiments on 7 common STS datasets:
438 STS tasks 2012-2016 (Agirre et al., 2012, 2013,
439 2014, 2015, 2016) STS-B(Cer et al., 2017), SICK-
440 R (Marelli et al., 2014). We use the SentEval
441 toolkit (Conneau and Kiela, 2018) to download all
442 7 datasets. The sentence pairs in each datasets are
443 scored from 0 to 5 to indicate semantic similarity.

444

5.2 Baselines

445 We compare our method with both enlightening and
446 state-of-the-art methods. To validate the effective-
447 ness of our method in the non fine-tuned setting, we
448 use the GLoVe (Pennington et al., 2014) and post-
449 process methods: BERT-flow (Li et al., 2020) and
450 BERT-whitening (Su et al., 2021) as baselines. For
451 the fine-tuned setting, we compare our method with
452 IS-BERT(Zhang et al., 2020), InferSent(Conneau
453 et al., 2017), Universal Sentence Encoder(Cer et al.,
454 2018), SBERT(Reimers and Gurevych, 2019) and
455 the contrastive learning based methods: SimCSE
456 (Gao et al., 2021) and ConSERT (Yan et al., 2021).

457

5.3 Implementation Details

458 For the non fine-tuned setting, we report the re-
459 sult of BERT to validate the effectiveness of our
460 representation method. For the fine-tuned setting,
461 we use BERT and RoBERTa with the same unsup-
462ervised and supervised training data with (Gao
463 et al., 2021). Our methods are trained with prompt
464 based contrastive learning with template denoising.
465 The templates used for both settings are manual
466 searched according to Table 4. More details can be
467 found in Appendix C.

468

5.4 Non Fine-Tuned BERT Results

469 To connect with the previous analysis of the poor
470 performance of original BERT, we report our
471 prompt based methods with non fine-tuned BERT
472 in Table 5. Using templates can substantially im-
473 prove the results of original BERT on all datasets.
474 Compared to pooling methods like averaging of last
475 layer or averaging of first and last layers, our meth-
476 ods can improve spearman correlation by more

Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
GLoVe embeddings avg. [†]	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
BERT last avg.	30.87	59.89	47.73	60.29	63.73	47.29	58.22	52.57
static avg.	42.38	56.74	50.60	65.08	62.39	56.82	58.15	56.02
first-last avg. [†]	39.70	59.38	49.67	66.03	66.19	53.87	62.06	56.70
static remove biases avg.	53.09	66.48	65.09	69.80	67.85	61.60	57.80	63.10
BERT-flow [†]	58.40	67.10	60.85	75.16	71.22	68.66	64.47	66.55
BERT-whitening [†]	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
Prompt based BERT (manual)	60.96	73.83	62.18	71.54	68.68	70.60	67.16	67.85
Prompt based BERT (manual&OptiPrompt)	64.56	79.96	70.05	79.37	75.35	77.25	68.56	73.59

Table 5: The performance comparison of our unfine-tuned BERT method on STS tasks. [†]: results from (Gao et al., 2021). The BERT-flow(Li et al., 2020) and BERT-whitening (Su et al., 2021) use the "NLI" setting. All BERT based methods use *bert-base-uncased*. Last avg. denotes averaging the last layer of BERT. Static avg. denotes averaging the static token embedding of BERT. First-last avg. (Su et al., 2021) uses the first and last layer. Static remove biases avg. means removing biased tokens in static avg., which we have introduced before.

Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
<i>Unsupervised models</i>								
IS-BERT _{base} [¶]	56.77	69.24	61.21	75.23	70.16	69.21	64.25	66.58
ConSERT _{base} [‡]	64.64	78.49	69.07	79.72	75.95	73.97	67.31	72.74
SimCSE-BERT _{base} [‡]	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
PromptBERT _{base}	71.56 _{±0.18}	84.58 _{±0.22}	76.98 _{±0.26}	84.47 _{±0.24}	80.60 _{±0.21}	81.60 _{±0.22}	69.87 _{±0.40}	78.54 _{±0.15}
RoBERTa _{base} -whitening [†]	46.99	63.24	57.23	71.36	68.99	61.36	62.91	61.73
SimCSE-RoBERTa _{base} [†]	70.16	81.77	73.24	81.36	80.65	80.22	68.56	76.57
PromptRoBERTa _{base}	73.94 _{±0.90}	84.74 _{±0.36}	77.28 _{±0.41}	84.99 _{±0.25}	81.74 _{±0.29}	81.88 _{±0.37}	69.50 _{±0.57}	79.15 _{±0.25}
<i>Supervised models</i>								
InferSent-GloVe [§]	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
SBERT _{base} [§]	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT _{base} -flow [†]	69.78	77.27	74.35	82.01	77.46	79.12	76.21	76.60
SBERT _{base} -whitening [†]	69.65	77.57	74.66	82.27	78.39	79.52	76.91	77.00
ConSERT _{base} [‡]	74.07	83.93	77.05	83.66	78.76	81.36	76.77	79.37
SimCSE-BERT _{base} [†]	75.30	84.67	80.19	85.40	80.82	84.25	80.39	81.57
PromptBERT _{base}	75.48	85.59	80.57	85.99	81.08	84.56	80.52	81.97
SRoBERTa _{base} [§]	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa _{base} -whitening [†]	70.46	77.07	74.46	81.64	76.43	79.49	76.65	76.60
SimCSE-RoBERTa _{base} [†]	76.53	85.21	80.95	86.03	82.57	85.83	80.50	82.52
PromptRoBERTa _{base}	76.75	85.93	82.28	86.69	82.80	86.14	80.04	82.95

Table 6: The performance comparison of our fine-tuned BERT methods on STS tasks. For unsupervised models, we found the result of unsupervised contrastive learning is unstable, and we train our model with 10 random seeds. [†]: results from (Gao et al., 2021). [‡]: results from (Yan et al., 2021). [§]: results from (Reimers and Gurevych, 2019). [¶]: results from (Zhang et al., 2020).

than 10%. Compared to the postprocess methods: BERT-flow and BERT-whitening, only using the manual template surpasses can these methods. Moreover, we can use the continuous template by OptiPrompt to help original BERT achieve much better results, which even outperforms unsupervised ConSERT in Table 6.

5.5 Fine-Tuned BERT Results

The results of fine-tuned BERT are shown in Table 6. Following previous works (Reimers and Gurevych, 2019), we run unsupervised and supervised methods respectively. Although the current contrastive learning based methods (Gao et al.,

2021; Yan et al., 2021) achieved significant improvement compared to the previous methods, our method still outperforms them. Prompt based contrastive learning objective significantly shortens the gap between the unsupervised and supervised methods. It also proves our method can leverage the knowledge of unlabeled data with different templates as positive pairs. Moreover, we report the unsupervised performance with 10 random seeds to achieve more accurate results. In Discussion, we also report the result of SimCSE with 10 random seeds. Compared to SimCSE, our method shows more stable results than it.

Sentence	Top-5 tokens	Top-5 tokens after template denoising
i am sad.	sad,sadness,happy,love,happiness	sad,sadness,crying,grief,tears
i am not happy.	happy,happiness,sad,love,nothing	sad,happy,unhappy,upset,angry
the man is playing the guitar.	guitar,song,music,guitarist,bass	guitar,guitarist,guitars,playing,guitarists
the man is playing the piano.	piano,music,no,yes,bass	piano,pianist,pianos,playing,guitar

Table 7: The top-5 tokens predicted by manual template with original BERT.

5.6 Effectiveness of Prompt Based Contrastive 503 Learning with Template Denoising

504 We report the results of different unsupervised training
505 objectives in prompt based BERT. We use the
506 following training objectives: 1) the same template,
507 which uses inner dropout noise as data augmentation
508 (Gao et al., 2021) 2) the different templates as
509 positive pairs 3) the different templates with tem-
510 plate denoising (our default method). Moreover,
511 we use the same template and setting to predict and
512 only change the way to generate positive pairs in
513 the training stage. All results are from 10 random
514 runs. The result is shown in Table 8. We observe
515 our method can achieve the best and most stable
516 results among three training objectives.
517

	BERT _{base}	RoBERTa _{base}
same template (dropout)	78.16 _{±0.17}	78.16 _{±0.44}
different templates	78.19 _{±0.29}	78.17 _{±0.44}
different templates with denoising	78.54 _{±0.15}	79.15 _{±0.25}

518 Table 8: Comparison of different unsupervised training
519 objectives.

520 6 Discussion

521 6.1 Template Denoising

522 We find the template denoising efficiently removes
523 the bias from templates and improves the quality
524 of top-k tokens predicted by MLM head in original
525 BERT. As Table 7 shows, we predict some
526 sentences’ top-5 tokens in the [MASK] tokens. We
527 find the template denoising removes the unrelated
528 tokens like “nothing,no,yes” and helps the model
529 predict more related tokens. To quantify this, we
530 also represent the sentence from the Eq. 3 by using
531 the weighted average of top-200 tokens as the sen-
532 tence embeddings. The results are shown in Table 9.
533 The template denoising significantly improves the
534 quality of tokens predicted by MLM head. How-
535 ever, it can’t improve the performance for our de-
536 fault represent method in the Eq. 2 ([MASK] token
537 in Table 9). In this work, we only use the tem-
538 plate denoising in our contrastive training objective,

539 which helps us eliminate different template biases.

	no denoising	denoising
avg. Top-200 tokens	56.19	60.39
[MASK] token	67.85	67.43

540 Table 9: Influence of template denoising in sentence
541 embeddings.

542 6.2 Stability in Unsupervised Contrastive 543 Learning

544 To prove the unstable results in unsupervised con-
545 trastive learning in sentence embeddings, we also
546 reproduce the result of unsupervised SimCSE-
547 BERT_{base} with 10 random seeds in Table 10. Our
548 results are more stable than SimCSE. The differ-
549 ence between the best and worst results can be up to
550 3.14% in SimCSE. However, the gap in our method
551 is only 0.53.

	Mean	Max	Min
SimCSE-BERT _{base}	75.42 _{±0.86}	76.64	73.50
PromptBERT _{base}	78.54 _{±0.15}	78.86	78.33

552 Table 10: Results in unsupervised contrastive learning.

553 7 Conclusion

554 In this paper, we analyzed the poor performance of
555 original BERT for sentence embeddings. The main
556 reason is not the anisotropy, but the static token
557 embeddings biases and ineffectively using original
558 BERT layers. Based on these findings, we pro-
559 posed a prompt based sentence embedding method
560 to avoid static token embeddings biases and lever-
561 age the pre-trained knowledge in the original BERT
562 layers. Our method significantly improved the per-
563 formance of the original BERT. We also proposed
564 a new prompt based contrastive learning method
565 to shorten the gap between the unsupervised and
566 supervised methods. Both our unsupervised and
567 supervised methods achieve the state-of-the-art per-
568 formance.

References

- 564 Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, et al. 2018. Universal sentence encoder for english. 621
 565 Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Inigo Lopez-Gazpio, Montse Maritxalar, Rada Mi- 622
 566 halcea, et al. 2015. Semeval-2015 task 2: Semantic 623
 567 textual similarity, english, spanish and pilot on inter- 624
 568 pretability. In *Proceedings of the 9th international 625
 569 workshop on semantic evaluation (SemEval 2015)*, 626
 570 pages 252–263. 627
- 573 Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, 628
 574 Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, 629
 575 Rada Mihalcea, German Rigau, and Janyce Wiebe. 630
 576 2014. Semeval-2014 task 10: Multilingual semantic 631
 577 textual similarity. In *Proceedings of the 8th interna- 632
 578 tional workshop on semantic evaluation (SemEval 2014)*, 633
 579 pages 81–91. 634
- 580 Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, 635
 581 Aitor Gonzalez Agirre, Rada Mihalcea, German 636
 582 Rigau Claramunt, and Janyce Wiebe. 2016. Semeval- 637
 583 2016 task 1: Semantic textual similarity, monolin- 638
 584 gual and cross-lingual evaluation. In *SemEval-2016; 639
 585 10th International Workshop on Semantic Evalu- 640
 586 ation; 2016 Jun 16-17; San Diego, CA. Stroudsburg 641
 587 (PA): ACL; 2016. p. 497-511.* ACL (Association for 642
 588 Computational Linguistics). 643
- 589 Eneko Agirre, Daniel Cer, Mona Diab, and Aitor 644
 590 Gonzalez-Agirre. 2012. Semeval-2012 task 6: A 645
 591 pilot on semantic textual similarity. In **SEM 2012: 646
 592 The First Joint Conference on Lexical and Compu- 647
 593 tational Semantics—Volume 1: Proceedings of the 648
 594 main conference and the shared task, and Volume 649
 595 2: Proceedings of the Sixth International Workshop 650
 596 on Semantic Evaluation (SemEval 2012)*, pages 385– 651
 597 393. 652
- 598 Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez- 653
 599 Agirre, and Weiwei Guo. 2013. * sem 2013 shared 654
 600 task: Semantic textual similarity. In *Second joint 655
 601 conference on lexical and computational semantics 656
 602 (*SEM), volume 1: proceedings of the Main confer- 657
 603 ence and the shared task: semantic textual similarity*, 658
 604 pages 32–43. 659
- 605 Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2016. A 660
 606 simple but tough-to-beat baseline for sentence em- 661
 607 beddings. 662
- 608 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie 663
 609 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind 664
 610 Neelakantan, Pranav Shyam, Girish Sastry, Amanda 665
 611 Askell, et al. 2020. Language models are few-shot 666
 612 learners. *arXiv preprint arXiv:2005.14165*. 667
- 613 Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez- 668
 614 Gazpio, and Lucia Specia. 2017. Semeval-2017 669
 615 task 1: Semantic textual similarity-multilingual and 670
 616 cross-lingual focused evaluation. *arXiv preprint 671
 617 arXiv:1708.00055*. 672
- 618 Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, 673
 619 Nicole Limtiaco, Rhomni St John, Noah Constant, 674
 620 Mario Guajardo-Cespedes, Steve Yuan, Chris Tar,
- et al. 2018. Universal sentence encoder for english. 621
 In *Proceedings of the 2018 Conference on Empirical 622
 Methods in Natural Language Processing: System 623
 Demonstrations*, pages 169–174. 624
- Alexis Conneau and Douwe Kiela. 2018. Senteval: An 625
 evaluation toolkit for universal sentence representa- 626
 tions. *arXiv preprint arXiv:1803.05449*. 627
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic 628
 Barrault, and Antoine Bordes. 2017. Supervised 629
 learning of universal sentence representations from 630
 natural language inference data. *arXiv preprint 631
 arXiv:1705.02364*. 632
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 633
 Kristina Toutanova. 2018. Bert: Pre-training of deep 634
 bidirectional transformers for language understand- 635
 ing. *arXiv preprint arXiv:1810.04805*. 636
- Kawin Ethayarajh. 2019. How contextual are context- 637
 ualized word representations? comparing the geo- 638
 metry of bert, elmo, and gpt-2 embeddings. *arXiv 639
 preprint arXiv:1909.00512*. 640
- Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie- 641
 Yan Liu. 2019. Representation degeneration problem 642
 in training natural language generation models. *arXiv 643
 preprint arXiv:1907.12009*. 644
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. 645
 Making pre-trained language models better few-shot 646
 learners. *arXiv preprint arXiv:2012.15723*. 647
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. 648
 Simcse: Simple contrastive learning of sentence em- 649
 beddings. *arXiv preprint arXiv:2104.08821*. 650
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, 651
 Yiming Yang, and Lei Li. 2020. On the sentence 652
 embeddings from bert for semantic textual similarity. 653
 In *Proceedings of the 2020 Conference on Empirical 654
 Methods in Natural Language Processing (EMNLP)*, 655
 pages 9119–9130. 656
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man- 657
 dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, 658
 Luke Zettlemoyer, and Veselin Stoyanov. 2019. 659
 Roberta: A robustly optimized bert pretraining 660
 approach. *arXiv preprint arXiv:1907.11692*. 661
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa 662
 Bentivogli, Raffaella Bernardi, Roberto Zamparelli, 663
 et al. 2014. A sick cure for the evaluation of com- 664
 positional distributional semantic models. In *Lrec*, 665
 pages 216–223. Reykjavik. 666
- Jeffrey Pennington, Richard Socher, and Christopher D 667
 Manning. 2014. Glove: Global vectors for word 668
 representation. In *Proceedings of the 2014 conference 669
 on empirical methods in natural language processing 670
 (EMNLP)*, pages 1532–1543. 671
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: 672
 Sentence embeddings using siamese bert-networks. 673
arXiv preprint arXiv:1908.10084. 674

675 Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou.
676 2021. Whitening sentence representations for bet-
677 ter semantics and faster retrieval. *arXiv preprint*
678 *arXiv:2103.15316*.

679 Hayato Tsukagoshi, Ryohei Sasano, and Koichi Takeda.
680 2021. Defsent: Sentence embeddings using defini-
681 tion sentences. *arXiv preprint arXiv:2105.04339*.

682 Alex Wang, Amanpreet Singh, Julian Michael, Felix
683 Hill, Omer Levy, and Samuel R Bowman. 2018.
684 Glue: A multi-task benchmark and analysis platform
685 for natural language understanding. *arXiv preprint*
686 *arXiv:1804.07461*.

687 Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le,
688 Mohammad Norouzi, Wolfgang Macherey, Maxim
689 Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al.
690 2016. Google’s neural machine translation system:
691 Bridging the gap between human and machine trans-
692 lation. *arXiv preprint arXiv:1609.08144*.

693 Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang,
694 Wei Wu, and Weiran Xu. 2021. Consert: A con-
695 trastive framework for self-supervised sentence repre-
696 sentation transfer. *arXiv preprint arXiv:2105.11741*.

697 Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim,
698 and Lidong Bing. 2020. An unsupervised sentence
699 embedding method by mutual information maximiza-
700 tion. *arXiv preprint arXiv:2009.12061*.

701 Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021.
702 Factual probing is [mask]: Learning vs. learning to
703 recall. *arXiv preprint arXiv:2104.05240*.

704 A Static Token Embeddings Biases

705 A.1 Eliminating Biases by Removing Tokens

706 We reported the detailed implementation of elimi-
707 nating static token embeddings biases by deleting
708 tokens on *bert-base-uncased*, *bert-base-cased* and
709 *roberta-base*. For Freq. tokens, we follow the set-
710 tings in (Yan et al., 2021) and remove the top 36
711 frequent tokens. The removed Freq. tokens are
712 shown in Table 11. For Sub. tokens, we directly
713 remove all subword tokens (yellow tokens in Fig-
714 ure 2). For Case. tokens, only SICK(Marelli et al.,
715 2014) has sentences with upper and lower case,
716 and we lowercase these sentences to remove the
717 uppercased tokens (red tokens in Figure 2). For
718 Pun., we remove the tokens, which contain only
719 punctuations.

720 A.2 Eliminating Biases by Pre-training

721 According to (Gao et al., 2019), we find the most
722 of biases in static token embeddings are gradient
723 from the MLM classification head weight, which
724 transform the last hidden vector of [MASK] to
725 the probability of all tokens. The tying weight

Removed Top frequency Tokens	
<i>bert-base-uncased</i>	. a the in , is to of and ' on and - s with for " at s woman are two that you dog said playing an as was from : by white
<i>bert-base-cased</i>	Ġ. Ġa Ġthe Ġin a Ġ, Ġis Ġto Ġof Ġon Ġ' s . the Ġman - Ġwith Ġfor Ġwoman Ġare Ġ" Ġthat Ġit Ġdog Ġplaying Ġwas Ġas Ġfrom Ġ: Ġyou i Ġby
<i>roberta-base</i>	

Table 11: Removed top 36 frequent tokens in *bert-base-cased*, *bert-base-uncased* and *roberta-base*.

726 between the static token embeddings and MLM
727 classification head causes static token embeddings
728 to suffer from bias problems.

729 We have pre-trained two BERT-like models with
730 the MLM pre-training objective. The only differ-
731 ence between the two pre-trained models is tying
732 and untying the weight between static token em-
733 beddings and MLM classification head. We have
734 pre-trained these two models on 125k steps with
735 2k batch sizes.

736 As shown in Figure 2, we have shown the static
737 token embeddings of the untying model, MLM
738 head weight of untying model and static token em-
739 beddings (MLM head weight) of the tying model.
740 The distribution of the tying model and the head
741 weight of the untying model is same with *bert-base-
742 cased* in Figure 1, which severely suffers from the
743 embedding biases. However, the distribution of the
744 token embeddings in the untying weights model is
745 less influenced by these biased. We also report the
746 average spearman correlation of three embedding
747 on STS tasks in Table 12. **Static token embeddings
748 of the untying model achieves the best correlation
749 among the three embedding.**

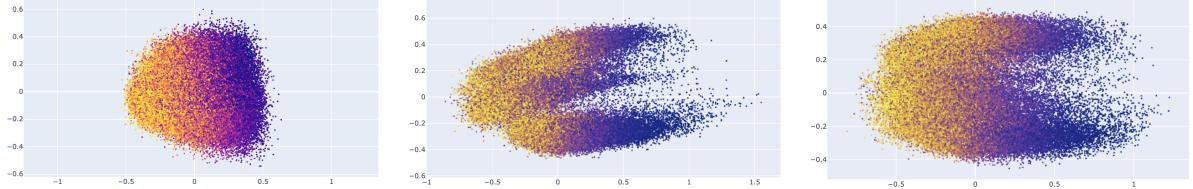
	Avg.
MLM head of untying model	43.33
Static token embeddings of untying model	49.41
Static token embeddings of tying model	45.68

Table 12: The avg. spearman correlation of three em-
beddings.

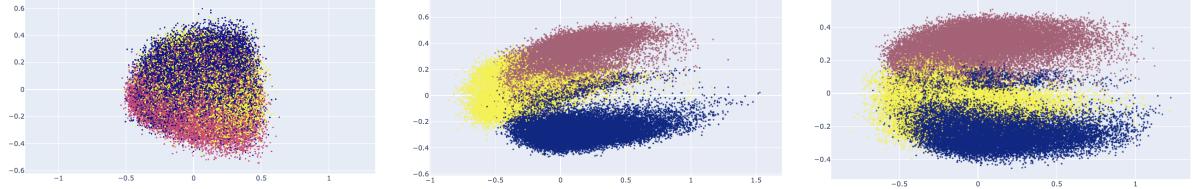
750 B Transfer Tasks

751 We also evaluate our models on the following trans-
752 fer tasks: MR, CR, SUBJ, MPQA, SST-2, TREC
753 and MRPC. We follow the default configurations in
754 SentEval⁴. The results is shown in Table 13. Com-
755 paring to SimCSE, our RoBERTa based method

⁴<https://github.com/facebookresearch/SentEval>



(a) Frequency bias in static token embeddings of untying weights pre-trained model.
(b) Frequency bias in MLM head of untying weights pre-trained model.
(c) Frequency bias in tying weights pre-trained model.



(d) Subword and Case biases in static token embeddings of untying weights pre-trained model.
(e) Subword and Case biases in MLM head of untying weights pre-trained model.
(f) Subword and Case biases in tying weights pre-trained model.

Figure 2: 2D visualization of static token embeddings in untying and tying weights pre-trained model. For frequency bias, the darker the color, the higher the token frequency. For subword and case bias, yellow represents subword and red represents the token contains capital letters.

Method	MR	CR	SUBJ	MPQA	SST-2	TREC	MRPC	Avg.
<i>Unsupervised models</i>								
Avg. BERT embedding	78.66	86.25	94.37	88.66	84.40	92.80	69.54	84.94
BERT-[CLS] embedding	78.68	84.85	94.21	88.23	84.13	91.40	71.13	84.66
IS-BERT	81.09	87.18	94.96	88.75	85.96	88.64	74.24	85.83
SimCSE-BERT	81.18	86.46	94.45	88.88	85.50	89.80	74.43	85.81
PromptBERT	80.74	85.49	93.65	89.32	84.95	88.20	76.06	85.49
SimCSE-RoBERTa	81.04	87.74	93.28	86.94	86.60	84.60	73.68	84.84
PromptRoBERTa	83.82	88.72	93.19	90.36	88.08	90.60	76.75	87.36
<i>Supervised models</i>								
InferSent-GloVe	81.57	86.54	92.50	90.38	84.18	88.20	75.77	85.59
Universal Sentence Encoder	80.09	85.19	93.98	86.70	86.38	93.20	70.14	85.10
SBERT	83.64	89.43	94.39	89.86	88.96	89.60	76.00	87.41
SimCSE-BERT	82.69	89.25	94.81	89.59	87.31	88.40	73.51	86.51
PromptBERT	83.14	89.38	94.49	89.93	87.37	87.40	76.58	86.90
SRoBERTa	84.91	90.83	92.56	88.75	90.50	88.60	78.14	87.76
SimCSE-RoBERTa	84.92	92.00	94.11	89.82	91.27	88.80	75.65	88.08
PromptRoBERTa	85.74	91.47	94.81	90.93	92.53	90.40	77.10	89.00

Table 13: Transfer task results of different sentence embedding models.

can improve 2.52 and 0.92 on unsupervised and supervised models respectively.

C Training Details

For the non fine-tuned setting, the manual template we used is *This sentence : “[X]” means [MASK]* .. For OptPrompt, we first initialize the template embeddings with the manual template and then train these template embeddings by freezing BERT

with the unsupervised training task followed by (Gao et al., 2021), and the batch size, learning-rate, epoch and valid steps are 256, 3e-5, 5 and 1000.

For the fine-tuned setting, all training data is same with (Gao et al., 2021). The max sentence sequence length is set to 32. For templates, we only use the manual templates, which are manually searched according to STS-B dev in unfine-tuned models. The templates is shown in Table 14. For unsupervised method, we use two different

756
757

758

759
760
761
762
763

764
765
766
767
768
769
770
771
772
773

774 templates for unsupervised training with template
 775 denoising according to our prompt based training
 776 objective. In predicting, we directly use the one
 777 template without template denoising . For super-
 778 vised method, we use template denoising with same
 779 template for contrastive learning, because we al-
 780 ready have supervised negative samples. We also
 781 report other training details in Table 15.

Model	Template
BERT	This sentence of “[X]” means [MASK] . This sentence : “[X]” means [MASK] .
RoBERTa	This sentence : ‘[X]’ means [MASK] . The sentence : ‘[X]’ means [MASK] .

Table 14: Templates for our method in fine-tuned setting

	<i>Unsupervised</i>		<i>Supervised</i>	
	BERT	RoBERTa	BERT	RoBERTa
Batch size	256	256	512	512
Learning rate	1e-5	1e-5	5e-5	5e-5
Epoch	1	1	3	3
Vaild steps	125	125	125	125

Table 15: Hyperparameters for our method in fine-tuned setting